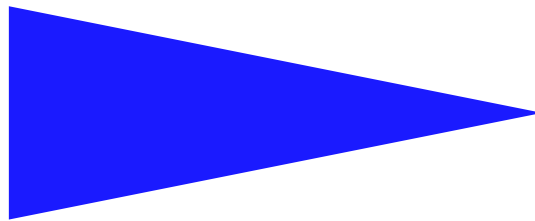


PUBLICATION  
INTERNE  
N° 1821



**ELECTING AN EVENTUAL LEADER  
IN AN ASYNCHRONOUS SHARED MEMORY SYSTEM**

**A. FERNÁNDEZ E. JIMÉNEZ M. RAYNAL**



# Electing an Eventual Leader in an Asynchronous Shared Memory System

A. Fernández\* E. Jiménez\*\* M. Raynal\*\*\*

Systèmes communicants

Publication interne n° 1821 — Novembre 2006 — 19 pages

**Abstract:** This paper considers the problem of electing an eventual leader in an asynchronous shared memory system. While this problem has received a lot of attention in message-passing systems, very few solutions have been proposed for shared memory systems. As an eventual leader cannot be elected in a pure asynchronous system prone to process crashes, the paper first proposes to enrich the asynchronous system model with an additional assumption. That assumption, denoted *AWB*, requires that after some time (1) there is a process whose write accesses to some shared variables are timely, and (2) the timers of the other processes are asymptotically well-behaved. The *asymptotically well-behaved* timer notion is a new notion that generalizes and weakens the traditional notion of timers whose durations are required to monotonically increase when the values they are set to increase. Then, the paper presents two *AWB*-based algorithms that elect an eventual leader. Both algorithms are independent of the value of  $t$  (the maximal number of processes that may crash). The first algorithm enjoys the following noteworthy properties: after some time only the elected leader has to write the shared memory, and all but one shared variables have a bounded domain, be the execution finite or infinite. This algorithm is consequently optimal with respect to the number of processes that have to write the shared memory. The second algorithm enjoys the following property: all the shared variables have a bounded domain. This is obtained at the following additional price: all the processes are required to forever write the shared memory. A theorem is proved that states this price has to be paid by any algorithm that elects an eventual leader in a bounded shared memory model. This second algorithm is consequently optimal with respect to the number of processes that have to write in such a constrained memory model. In a very interesting way, these algorithms show an inherent tradeoff relating the number of processes that have to write the shared memory and the bounded/unbounded attribute of that memory.

**Key-words:** Access cost, Asynchronous system, Bounded memory, Eventual leader, Fault-tolerance, Omega, Process crash, Shared memory, System model, Timer.

(Résumé : *tsvp*)

\* LADyR, GSyC, Universidad Rey Juan Carlos, 28933 Móstoles, Spain, anto@gsyc.escet.urjc.es

The work of this author was done while on leave at IRISA, supported by the Spanish MEC under grant PR-2006-0193.

\*\* EUI, Universidad Politécnica de Madrid, 28031 Madrid, Spain, ernes@eui.upm.es.

The work of A. Fernández and E. Jiménez was partially supported by the Spanish MEC under grants TIN2005-09198-C02-01, TIN2004-07474-C02-02, and TIN2004-07474-C02-01, and the Comunidad de Madrid under grant S-0505/TIC/0285.

\*\*\* IRISA, Université de Rennes 1, Campus de Beaulieu, 35042, Rennes Cedex, France, raynal@irisa.fr



# **Election d'un leader dans un système à mémoire partagée asynchrone**

**Résumé :** Ce rapport présente un protocole d'élection d'un leader inéluctable dans un système asynchrone à mémoire partagée dans lequel un nombre quelconque de processus peuvent crasher.

**Mots clés :** Systèmes asynchrones, Tolérance aux fautes, Crash de processus, Oracle oméga, Mémoire partagée, Leader inéluctable, Chien de garde.

# 1 Introduction

**Equipping an asynchronous system with an oracle** An asynchronous system is characterized by the absence of a bound on the time it takes for a process to proceed from a step of its algorithm to the next one. Combined with process failures, such an absence of a bound can make some synchronization or coordination problems impossible to solve (even when the processes communicate through a reliable communication medium). The most famous of these “impossible” asynchronous problems is the well-known *consensus* problem [8]. Intuitively, this impossibility comes from the fact that a process cannot safely distinguish a crashed process from a very slow process.

One way to address and circumvent these impossibilities consists on enriching the underlying asynchronous systems with an appropriate *oracle* [26]. More precisely, in a system prone to process failures, such an oracle (sometimes called *failure detector*) provides each process with hints on which processes are (or are not) faulty. According to the quality of these hints, several classes of oracles can be defined [5]. So, given an asynchronous system prone to process failures equipped with an appropriate oracle, it becomes possible to solve a problem that is, otherwise, impossible to solve in a purely asynchronous system. This means that an oracle provides processes with additional computability power.

**Fundamental issues related to oracles for asynchronous systems** Two fundamental questions can be associated with oracles. The first is more on the theoretical side and concerns their computability power. Given a problem (or a family of related problems), which is the weakest oracle that allows solving that problem in an asynchronous system where processes can experience a given type of failures? Intuitively, an oracle  $O_w$  is the weakest for solving a problem  $P$  if it allows solving that problem, and any other oracle  $O_{nw}$  that allows solving  $P$  provides hints on failures that are at least as accurate as the ones provided by  $O_w$  (this means that the properties defining  $O_{nw}$  imply the ones defining  $O_w$ , but not necessarily vice-versa). It has been shown that, in asynchronous systems prone to process crash failures, the class of *eventual leader* oracles is the weakest for solving asynchronous *consensus*, be these systems message-passing systems [6] or shared memory systems [19]<sup>1</sup>. It has also been shown that, for the same type of process failures, the class of *perfect failure detectors* (defined in [5]) is the weakest for solving asynchronous *interactive consistency* [14].

The second important question is on the algorithm/protocol side and concerns the implementation of oracles (failure detectors) that are designed to equip an asynchronous system. Let us first observe that no such oracle can be implemented on top of a purely asynchronous system (otherwise the problem it allows solving could be solved in a purely asynchronous system without additional computability power). So, this fundamental question translates as follows. First, find “reasonably weak” behavioral assumptions that, when satisfied by the underlying asynchronous system, allow implementing the oracle. “Reasonably weak” means that, although they cannot be satisfied by all the runs, the assumptions are actually satisfied in “nearly all” the runs of the asynchronous system. Second, once such assumptions have been stated, design efficient algorithms that implement correctly the oracle in all the runs satisfying the assumptions.

**Content of the paper** Considering the asynchronous shared memory model where any number of processes can crash, this paper addresses the construction of eventual leader oracles [6]. Such an oracle (usually denoted  $\Omega$ )<sup>2</sup> provides the processes with a primitive `leader()` that returns a process identity, and satisfies the following “eventual” property in each run  $R$ : There is a time after which all the invocations of `leader()` return the same identity, that is the identity of a process that does not crash in the run  $R$ .

As already indicated, such an oracle is the weakest to solve the consensus problem in an asynchronous system where processes communicate through single-writer/multi-readers (1WnR) atomic registers and are prone to crash failures [19].

The paper has two main contributions.

- It first proposes a behavioral assumption that is particularly weak. This assumption is the following one. In each run, there are a finite (but unknown) time  $\tau$  and a process  $p$  (not a priori known) that does not crash in that run, such that after  $\tau$ :

---

<sup>1</sup>Let us also notice that the Paxos fault-tolerant state machine replication algorithm [16] is based on the  $\Omega$  abstraction. For the interested reader, an introduction to the family of Paxos algorithms can be found in [12].

<sup>2</sup>Without ambiguity and according to the context,  $\Omega$  is used to denote either the class of eventual leader oracles, or an oracle of that class.

- (1) There is a bound  $\Delta$  (not necessarily known) such that any two consecutive accesses to some shared variables issued by  $p$  are separated by at most  $\Delta$  time units, and
- (2) Each correct process  $q \neq p$  has a timer that is *asymptotically well-behaved*. Intuitively, this notion expresses the fact that eventually the duration that elapses before a timer expires has to increase when the timeout parameter increases.

It is important to see that the timers can behave arbitrarily during arbitrarily long (but finite) periods. Moreover, as we will see in the formal definition, their duration are not required to strictly increase according to their timeout periods. After some time, they have only to be lower-bounded by some monotonously increasing function.

It is noteworthy to notice that no process (but  $p$ ) is required to have any synchronous behavior. Only their timers have to eventually satisfy some (weak) behavioral property.

- The paper then presents two algorithms that construct an  $\Omega$  oracle in all the runs that satisfy the previous behavioral assumptions, and associated lower bounds. All the algorithms use atomic 1WnR atomic registers. The algorithms, that are of increasing difficulty, are presented incrementally.
  - In the first algorithm, all (but one of) the shared variables have a bounded domain (the size of which depends on the run). More specifically, this means that, be the execution finite or infinite, even the timeout values stop increasing forever.
 

Moreover, after some time, there is a single process that writes the shared memory. The algorithm is consequently write-efficient. It is even write-optimal as at least one process has to write the shared memory to inform the other processes that the current leader is still alive.
  - The second algorithm improves the first one in the sense that all the (local and shared) variables are bounded. This nice property is obtained by using two boolean flags for each pair of processes. These flags allow each process  $p$  to inform each other process  $q$  that it has read some value written by  $q$ .
  - Lower bound results are proved for the considered model. Two theorems are proved that state (1) the process that is eventually elected has to forever write the shared memory, and (2) any process (but the eventual leader) has to forever read from the shared memory. Another theorem shows that, if the shared memory is bounded, then all the processes have to forever write into the shared memory. There theorems show that both the algorithms presented in the paper are optimal with respect to these criteria.

**Why shared memory-based  $\Omega$  algorithms are important** Some distributed systems are made up of computers that communicate through a network of attached disks. These disks constitute a storage area network (SAN) that implements a shared memory abstraction. As commodity disks are cheaper than computers, such architectures are becoming more and more attractive for achieving fault-tolerance. The  $\Omega$  algorithms presented in this paper are suited to such systems [1, 4, 10, 18].

**Related work** As far as we know, a single shared memory  $\Omega$  algorithm has been proposed so far [13]. This algorithm considers that the underlying system satisfies the following behavioral assumption: there is a time  $\tau$  after which there are a lower bound and an upper bound for any process to execute a local step, or a shared memory access. This assumption defines an eventually synchronous shared memory system. It is easy to see that thus is a stronger assumption than the assumption previously defined here.

The implementation of  $\Omega$  in asynchronous message-passing systems is an active research area. Two main approaches have been investigated: the *timer*-based approach and the *message pattern*-based approach.

The timer-based approach relies on the addition of timing assumptions [7]. Basically, it assumes that there are bounds on process speeds and message transfer delays, but these bounds are not known and hold only after some finite but unknown time. The algorithms implementing  $\Omega$  in such “augmented” asynchronous systems are based on timeouts (e.g., [2, 3, 17]). They use successive approximations to eventually provide each process with an upper bound on transfer delays and processing speed. They differ mainly on the “quantity” of additional synchrony they consider, and on the message cost they require after a leader has been elected.

Among the protocols based on this approach, a protocol presented in [2] is particularly attractive, as it considers a relatively weak additional synchrony requirement. Let  $t$  be an upper bound on the number of processes that may crash ( $1 \leq t < n$ , where  $n$  is the total number of processes). This assumption is the following: the underlying asynchronous system, which can have fair lossy channels, is required to have a correct process  $p$  that is a  $\diamond t$ -source. This means that  $p$  has  $t$  output channels that are eventually timely: there is a time after which the transfer delays of all the messages sent on such a channel are bounded (let us notice that this is trivially satisfied if the receiver has crashed). Notice that such a  $\diamond t$ -source is not known in advance and may never be explicitly known. It is also shown in [2] that there is no leader protocol if the system has only  $\diamond(t-1)$ -sources. A versatile adaptive timer-based approach has been developed in [20].

The message pattern-based approach, introduced in [21], does not assume eventual bounds on process and communication delays. It considers that there is a correct process  $p$  and a set  $Q$  of  $t$  processes (with  $p \notin Q$ , moreover  $Q$  can contain crashed processes) such that, each time a process  $q \in Q$  broadcasts a query, it receives a response from  $p$  among the first  $(n-t)$  corresponding responses (such a response is called a winning response). It is easy to see that this assumption does not prevent message delays to always increase without bound. Hence, it is incomparable with the synchrony-related  $\diamond t$ -source assumption. This approach has been applied to the construction of an  $\Omega$  algorithm in [23].

A *hybrid* algorithm that combines both types of assumption is developed in [24]. More precisely, this algorithm considers that each channel eventually is timely or satisfies the message pattern, without knowing in advance which assumption it will satisfy during a particular run. The aim of this approach is to increase the assumption coverage, thereby improving fault-tolerance [25].

**Roadmap** The paper is made up of 5 sections. Section 2 presents the system model and the additional behavioral assumption. Then, Sections 3 and 4 present in an incremental way the two algorithms implementing an  $\Omega$  oracle, and show they are optimal with respect to the number of processes that have to write or read the shared memory. Finally, Section 5 provides concluding remarks.

## 2 Base Model, Eventual Leader and Additional Behavioral Assumption

### 2.1 Base asynchronous shared memory model

The system consists of  $n$ ,  $n > 1$ , processes denoted  $p_1, \dots, p_n$ . The integer  $i$  denotes the identity of  $p_i$ . (Sometimes a process is also denoted  $p$ ,  $q$  or  $r$ .) A process can fail by *crashing*, i.e., prematurely halting. Until it possibly crashes, a process behaves according to its specification, namely, it executes a sequence of steps as defined by its algorithm. After it has crashed, a process executes no more steps. By definition, a process is *faulty* during a run if it crashes during that run; otherwise it is *correct* in that run. There is no assumption on the maximum number  $t$  of processes that may crash, which means that up to  $n-1$  process may crash in a run.

The processes communicate by reading and writing a memory made up of atomic registers (also called shared variables in the following). Each register is one-writer/multi-reader (1WnR). “1WnR” means that a single process can write into it, but all the processes can read it. (Let us observe that using 1WnR atomic registers is particularly suited for cached-based distributed shared memory.) The only process allowed to write an atomic register is called its owner. *Atomic* means that, although read and write operations on the same register may overlap, each (read or write) operation appears to take effect instantaneously at some point of the time line between its invocation and return events (this is called the *linearization* point of the operation) [15]. Uppercase letters are used for the identifiers of the shared registers. These registers are structured into arrays. As an example,  $PROGRESS[i]$  denotes a shared register that can be written only by  $p_i$ , and read by any process.

Some shared registers are *critical*, while other shared registers are not. A critical register is a an atomic register on which some constraint can be imposed by the additional assumptions that allow implementing an eventual leader. This attribute allows restricting the set of registers involved in these assumptions.

A process can have local variables. They are denoted with lowercase letters, with the process identity appearing as a subscript. As an example,  $candidates_i$  denotes a local variable of  $p_i$ .

This base model is characterized by the fact that there is no assumption on the execution speed of one process with respect to another. This is the classical *asynchronous* crash prone shared memory model. It is denoted  $\mathcal{AS}_n[\emptyset]$  in the following

## 2.2 Eventual leader service

The notion of *eventual leader* oracle has been informally presented in the introduction. It is an entity that provides each process with a primitive `leader()` that returns a process identity each time it is invoked. A unique correct leader is eventually elected but there is no knowledge of when the leader is elected. Several leaders can coexist during an arbitrarily long period of time, and there is no way for the processes to learn when this “anarchy” period is over. The *leader* oracle, denoted  $\Omega$ , satisfies the following property [6]:

- **Validity:** The value returned by a `leader()` invocation is a process identity.
- **Eventual Leadership**<sup>3</sup>: There is a finite time and a correct process  $p_i$  such that, after that time, every `leader()` invocation returns  $i$ .
- **Termination:** Any `leader()` invocation issued by a correct process terminates.

The  $\Omega$  leader abstraction has been introduced and formally developed in [6] where it is shown to be the weakest, in terms of information about failures, to solve consensus in asynchronous systems prone to process crashes (assuming a majority of correct processes). Several  $\Omega$ -based consensus protocols have been proposed (e.g., [11, 16, 22] for message-passing systems, and [9] for shared memory systems)<sup>4</sup>.

## 2.3 Additional behavioral assumption

**Underlying intuition** As already indicated,  $\Omega$  cannot be implemented in pure asynchronous systems such as  $\mathcal{AS}_n[\emptyset]$ . So, we consider the system is no longer fully asynchronous: its runs satisfy the following assumption denoted *AWB* (for *asymptotically well-behaved*). The resulting system is consequently denoted  $\mathcal{AS}_n[AWB]$ .

Each process  $p_i$  is equipped with a timer denoted *timer<sub>i</sub>*. The intuition that underlies *AWB* is that, once a process  $p_\ell$  is defined as being the current leader, it should not to be demoted by a process  $p_i$  that believes  $p_\ell$  has crashed. To that end, constraints have to be defined on the behavior of both  $p_\ell$  and  $p_i$ . The constraint on  $p_\ell$  is to force it to “regularly” inform the other processes that it is still alive. The constraint on a process  $p_i$  is to prevent it to falsely suspect that  $p_\ell$  has crashed.

There are several ways to define runs satisfying the previous constraints. As an example, restricting the runs to be “eventually synchronous” would work but is much more constraining than what is necessary. The aim of the *AWB* additional assumption is to state constraints that are “as weak as possible”<sup>5</sup>. It appears that requiring the timers to be eventually monotonous is stronger than necessary (as we are about to see, this is a particular case of the *AWB* assumption). The *AWB* assumption is made up of two parts *AWB<sub>1</sub>* and *AWB<sub>2</sub>* that we present now. *AWB<sub>1</sub>* is on the existence of a process whose behavior has to satisfy a synchrony property. *AWB<sub>2</sub>* is on the timers of the other processes. *AWB<sub>1</sub>* and *AWB<sub>2</sub>* are “matching” properties.

**The assumption *AWB<sub>1</sub>*** The *AWB<sub>1</sub>* assumption requires that eventually a process does not behave in a fully asynchronous way. It is defined as follows.

*AWB<sub>1</sub>*: There are a time  $\tau_{01}$ , a bound  $\Delta$ , and a correct process  $p_\ell$  ( $\tau_{01}$ ,  $\Delta$  and  $p_\ell$  may be never explicitly known) such that, after  $\tau_{01}$ , any two consecutive accesses issued by  $p_\ell$  to (its own) critical registers, are completed in at most  $\Delta$  time units.

This property means that, after some arbitrary (but finite) time, the speed of  $p_\ell$  is lower-bounded, i.e., its behavior is partially synchronous (let us notice that, while there is a lower bound, no upper bound is required on the speed of  $p_\ell$ , except the fact that it is not  $+\infty$ ).

---

<sup>3</sup>This property refers to a notion of global time. This notion is not accessible to the processes.

<sup>4</sup>It is important to notice that, albeit it can be rewritten using  $\Omega$  (first introduced in 1992), the original version of Paxos, that dates back to 1989, was not explicitly defined with this formalism.

<sup>5</sup>Of course, the notion of “as weak as possible” has to be taken with its intuitive meaning. This means that, when we want to implement  $\Omega$  in a shared memory system, we know neither an assumption weaker than *AWB*, nor the answer to the question: Is *AWB* the weakest additional assumption?

**The assumption  $AWB_2$**  In order to define  $AWB_2$ , we first introduce a function  $f()$  with monotonicity properties that will be used to define an asymptotic behavior. That function takes two parameters, a time  $\tau$  and a duration  $x$ , and returns a duration. It is defined as follows. There are two (possibly unknown) bounded values  $x_f$  and  $\tau_f$  such that:

- (f1)  $\forall \tau_2, \tau_1 : \tau_2 \geq \tau_1 \geq \tau_f, \forall x_2, x_1 : x_2 \geq x_1 \geq x_f: f(\tau_2, x_2) \geq f(\tau_1, x_1)$ . (After some point,  $f()$  is not decreasing with respect to  $\tau$  and  $x$ ).
- (f2)  $\lim_{x \rightarrow +\infty} f(\tau_f, x) = +\infty$ . (Eventually,  $f()$  always increases<sup>6</sup>.)

We are now in order to define the notion of *asymptotically well-behaved* timer. Considering the timer  $timer_i$  of a process  $p_i$  and a run  $R$ , let  $\tau$  be a real time at which the timer is set to a value  $x$ , and  $\tau'$  be the finite real time at which that timer expires. Let  $T_R(\tau, x) = \tau' - \tau$ , for each  $x$  and  $\tau$ . Then timer  $timer_i$  is asymptotically well-behaved in a run  $R$ , if there is a function  $f_R()$ , as defined above, such that:

- (f3)  $\forall \tau : \tau \geq \tau_f, \forall x : x \geq x_f: f_R(\tau, x) \leq T_R(\tau, x)$ .

This constraint states the fact that, after some point, the function  $T_R()$  is always above the function  $f_R()$ . It is important to observe that, after  $(\tau_f, x_f)$ , the function  $T_R(\tau, x)$  is not required to be non-decreasing, it can increase and decrease. Its only requirement is to always dominate  $f_R()$ . (See Figure 1.)

$AWB_2$ : The timer of each correct process (except possibly  $p_\ell$ ) is asymptotically well-behaved.

When we consider  $AWB$ , it is important to notice that any process (but  $p_\ell$  constrained by a speed lower bound) can behave in a fully asynchronous way. Moreover, the local clocks used to implement the timers are required to be neither synchronized, nor accurate with respect to real-time.

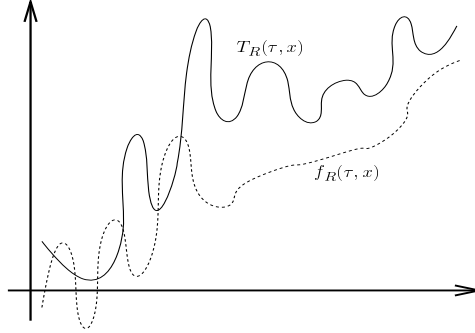


Figure 1:  $T_R()$  asymptotically dominates  $f_R()$

### 3 An $\Omega$ algorithm for $\mathcal{AS}_n[AWB]$

#### 3.1 Principles of the algorithm

The first algorithm implementing  $\Omega$  in  $\mathcal{AS}_n[AWB]$  that we present, relies on a very simple idea that has been used in several algorithms that build  $\Omega$  in message-passing systems. Each process  $p_i$  handles a set ( $candidates_i$ ) containing the processes that (from its point of view) are candidates for being the leader. When it suspects one of its candidates  $p_j$  to have crashed,  $p_i$  makes public the fact that it suspects  $p_j$  once more. (This is done by  $p_i$  increasing the shared register  $SUSPICIONS[i, j]$ .)

Finally, a process  $p_i$  defines its current leader as the least suspected process among its current candidates. As several processes can be equally suspected,  $p_i$  uses the function  $lexmin(X)$  that outputs the lexicographically smallest pair in the set parameter  $X$ , where  $X$  is the set of (number of suspicions, process identity) pairs defined from  $candidate_i$ , and  $(a, i) < (b, j)$  iff  $(a < b) \vee (a = b \wedge i < j)$ .

<sup>6</sup>If the image of  $f()$  is the set of natural numbers, then this condition can be replaced by  $x_2 > x_1 \implies f(\tau_f, x_2) > f(\tau_f, x_1)$ .

## 3.2 Description of the algorithm

The algorithm, based on the principles described just above, that builds  $\Omega$  in  $\mathcal{AS}_n[AWB]$  is depicted in Figure 2.

**Shared variables** The variables shared by the processes are the following:

- $SUSPICIONS[1..n, 1..n]$  is an array of natural registers.  $SUSPICIONS[j, k] = x$  means that, up to now,  $p_j$  has suspected  $x$  times the process  $p_k$  to have crashed. The entries  $SUSPICIONS[j, k]$ ,  $1 \leq k \leq n$  can be written only by  $p_j$ .
- $PROGRESS[1..n]$  is an array of natural registers. Only  $p_i$  can write  $PROGRESS[i]$ . (It does it only when it considers it is the leader.)
- $STOP[1..n]$  is an array of boolean registers. Only  $p_i$  can write  $STOP[i]$ . It sets it to *false* to indicate it considers itself as leader, and sets it to *true* to indicate it stops considering it is the leader.

The initial values of the previous shared variables can be arbitrary<sup>7</sup>. To improve efficiency, we consider that the natural integer variables are initialized to 0 and the boolean variables to *true*.

Each shared register  $PROGRESS[k]$  or  $STOP[k]$ ,  $1 \leq k \leq n$  is critical. Differently, none of the registers  $SUSPICIONS[j, k]$ ,  $1 \leq j, k \leq n$ , is critical. This means that, for a process  $p_k$  involved in the assumption  $AWB_1$ , only the accesses to its registers  $PROGRESS[k]$  and  $STOP[k]$  are concerned.

Let us observe that, as the shared variables  $PROGRESS[i]$ ,  $STOP[i]$  and  $SUSPICIONS[i, k]$ ,  $1 \leq k \leq n$ , are written only by  $p_i$ , that process can save their values in local memory and, when it has to read any of them, it can read instead its local copy. (We do not do it in our description of the algorithms to keep simpler the presentation.)

**Process behavior** The algorithm is made up of three tasks. Each local variable  $candidate_i$  is initialized to any set of process identities containing  $i$ .

The task  $T1$  implements the *leader()* primitive. As indicated,  $p_i$  determines the least suspected among the processes it considers as candidates (lines 2-4), and returns its identity (line 5).

The task  $T2$  is an infinite loop. When it considers it is the leader, (line 7),  $p_i$  repeatedly increases  $PROGRESS[i]$  to inform the other processes that it is still alive (lines 7-10). If it discovers it is no longer leader,  $p_i$  sets  $STOP[i]$  to *true* (line 11) to inform the other processes it is no longer competing to be leader.

Each process  $p_i$  has a local timer (denoted  $timer_i$ ), and manages a local variable  $last_i[k]$  where it saves the greatest value that it has ever read from  $PROGRESS[k]$ . The task  $T3$  is executed each time that timer expires (line 13). Then,  $p_i$  executes the following statements with respect to each process  $p_k$  (but itself, see line 14). First,  $p_i$  checks if  $p_k$  did some progress since the previous timer expiration (line 17). Then, it does the following.

- If  $PROGRESS[k]$  has progressed,  $p_i$  considers  $p_k$  as a candidate to be leader. To that end it adds  $k$  to the local set  $candidates_i$  (line 18). (It also updates  $last_i[k]$ , line 19.)
- If  $PROGRESS[k]$  has not progressed,  $p_i$  checks the value of  $STOP[k]$  (line 20). If it is true,  $p_k$  voluntarily demoted itself from being a candidate. Consequently,  $p_i$  suppresses  $k$  from its local set  $candidates_i$  (line 21). If  $STOP[k]$  is false and  $p_k$  is candidate from  $p_i$ 's point of view (line 22),  $p_i$  suspects  $p_k$  to have crashed (line 23) and suppresses it from  $candidates_i$  (line 24).

Then,  $p_i$  resets its local timer (line 27). Let us observe that no variable of the array  $SUSPICIONS$  can decrease and such an entry is increased each time a process is suspected by another process. Thanks to the properties, we will see in the proof that  $\max(\{SUSPICIONS[i, k]\}_{1 \leq k \leq n}) + 1$  can be used as the next timeout value. Note that to compute this value only variables owned by  $p_i$  are accessed.

## 3.3 Proof of the algorithm

**Lemma 1** *Let  $p_k$  be a faulty process and  $p_i$  a correct process. Eventually, the predicate  $k \notin candidates_i$  remains true forever.*

---

<sup>7</sup>This means that the algorithm is *self-stabilizing* with respect to the shared variables. Whatever their initial values, it converges in a finite number of steps towards a common leader, as soon as the additional assumption is satisfied.

```

task T1:
(1) when leader() is invoked:
(2)   for_each  $k \in candidates_i$  do
(3)      $susp_i[k] \leftarrow \Sigma_{1 \leq j \leq n} SUSPICIONS[j, k]$  end_for;
(4)   let  $(-, \ell) = \text{lex\_min}(\{susp_i[k], k\}_{k \in candidates_i});$ 
(5)   return( $\ell$ )

task T2:
(6) repeat_forever
(7)   while (leader() =  $i$ ) do
(8)      $PROGRESS[i] \leftarrow PROGRESS[i] + 1;$ 
(9)     if  $STOP[i]$  then  $STOP[i] \leftarrow \text{false}$  end_if
(10)  end_while;
(11) if  $(\neg STOP[i])$  then  $STOP[i] \leftarrow \text{true}$  end_if
(12) end_repeat

task T3:
(13) when  $timer_i$  expires:
(14) for_each  $k \in \{1, \dots, n\} \setminus \{i\}$  do
(15)    $stop\_k_i \leftarrow STOP[k];$ 
(16)    $progress\_k_i \leftarrow PROGRESS[k];$ 
(17)   if  $(progress\_k_i \neq last_i[k])$  then
(18)      $candidates_i \leftarrow candidates_i \cup \{k\};$ 
(19)      $last_i[k] \leftarrow progress\_k_i$ 
(20)   else_if  $(stop\_k_i)$  then
(21)      $candidates_i \leftarrow candidates_i \setminus \{k\}$ 
(22)   else_if  $(k \in candidates_i)$  then
(23)      $SUSPICIONS[i, k] \leftarrow SUSPICIONS[i, k] + 1;$ 
(24)      $candidates_i \leftarrow candidates_i \setminus \{k\}$ 
(25)   end_if
(26) end_for;
(27) set  $timer_i$  to  $\max(\{SUSPICIONS[i, k]\}_{1 \leq k \leq n}) + 1$ 

```

Figure 2: Write-efficient, all variables are 1WMR, bounded except a single entry of *PROGRESS*

**Proof** Let us consider a time  $\tau$  at which  $p_k$  has crashed. After  $\tau$ ,  $p_k$  never increases  $PROGRESS[k]$ . So, there is a time  $\tau' \geq \tau$  after which the test  $progress\_k_i \neq last_i[k]$  (line 17) is always false. It then follows from the lines 20-24 that, if  $k$  was in  $candidates_i$ , it is suppressed from this set. Moreover, as from now on we have  $last_i[k] = PROGRESS[k]$  forever, it follows from line 17 that  $k$  can never be again added to  $candidates_i$ .  $\square_{\text{Lemma 1}}$

Given a run  $R$  and a process  $p_x$ , let  $M_x$  denote the largest value of  $\Sigma_{1 \leq j \leq n} SUSPICIONS[j, x]$ . If there is no such value (i.e.,  $\Sigma_{1 \leq j \leq n} SUSPICIONS[j, x]$  grows forever), let  $M_x = +\infty$ . Finally, let  $B$  be the set of correct processes  $p_x$  such that  $M_x \neq +\infty$  ( $B$  stands for “bounded”).

**Lemma 2** *Let us assume that the behavioral assumption AWB is satisfied. Then,  $B \neq \emptyset$ .*

**Proof** Let  $p_i$  be a process that satisfies assumption  $AWB_i$ . (Hence,  $p_i$  is a correct process.) We show  $i \in B$ .

Let us first observe that the task *T3* of a process never executes the body of the **for** loop (lines 14-26) for that process, from which it follows that  $SUSPICIONS[i, i]$  is never increased. Let us now consider a faulty process  $p_j$ . Clearly, after  $p_j$  crashes,  $SUSPICIONS[j, i]$  is no longer increased.

So, the rest of the proof consists in showing that  $SUSPICIONS[j, i]$  remains bounded for any correct process  $p_j$ . Let  $S$  be the sequence of write operations issued by  $p_i$  at the lines 8 (write into  $PROGRESS[i]$ ), and 11 (write into  $STOP[i]$ )<sup>8</sup>. We consider two cases.

- $S$  is finite.

As  $p_i$  is correct, its last write operation in the sequence  $S$  is done at line 11, and that operation writes *true* into  $STOP[i]$ . Let  $\tau$  be the time at which  $p_i$  issued this last write operation. Since, after  $\tau$ ,  $PROGRESS[i]$  is no

<sup>8</sup>Let us notice that the write operations into  $STOP[i]$  at line 9 do not appear in the definition of  $S$ .

longer increased and  $STOP[i]$  is always equal to *true*, it follows that there is a time  $\tau' \geq \tau$  after which, when  $p_j$  executes line 17 with  $k = i$ , the test is always negative, while when it executes 20, the test is always positive. It follows that after  $\tau'$ ,  $p_j$  never executes line 23 with  $k = i$ , from which we conclude that  $SUSPICIONS[j, i]$  is bounded.

- $S$  is infinite.

Due to the assumption  $AWB_1$ , there are a time  $\tau_{01}$  and a bound  $\Delta$  such that, after  $\tau_{01}$ , any two consecutive operations on its critical registers issued by  $p_i$  are completed in at most  $\Delta$  time units. Observe that, after  $\tau_{01}$ , two updates of  $PROGRESS[i]$  (line 8) that are adjacent in  $S$  are completed in at most  $3\Delta$  time units<sup>9</sup>. Similarly, after  $\tau_{01}$ , an update of  $PROGRESS[i]$  (line 8) followed in  $S$  by an update of  $STOP[i]$  at line 11, is completed in at most  $3\Delta$  time units.

Figure 3 shows a possible execution after  $\tau_{01}$ . L8 and L11 represents consecutive updates of the 1WnR atomic registers  $PROGRESS[i]$  and  $STOP[i]$  at line 8 and line 11, respectively. The update of  $STOP[i]$  at line 9 is indicated with a black dot, and the area with stripes shows the interval during which the value of  $STOP[i]$  is *false*.

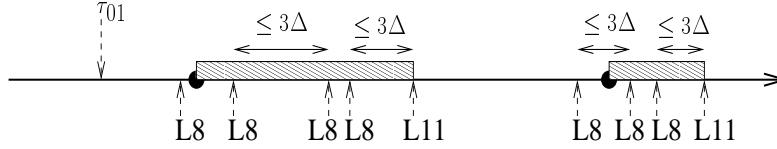


Figure 3: Illustrating the sequence  $S$

By assumption  $AWB_2$  we have that the timer of  $p_j$  is asymptotically well-behaved and for each run  $R$  there is a corresponding function  $f_R()$ , and parameters  $x_f$  and  $\tau_f$ . Let  $x_0 \geq x_f$  be a finite value such that  $f_R(\tau_f, x_0) = \Delta' > 3\Delta$ . Assumption (f2) implies that such a value of  $x_0$  always exists.

If  $SUSPICIONS[j, i]$  never reaches  $x_0$ , it is bounded and the lemma follows. So, let us consider that  $SUSPICIONS[j, i]$  reaches  $x_0$ . There is then a time after which  $timer_j$  is always set to a value greater than  $x_0$  (line 27). Let  $\tau_1 \geq \tau_f$  be a time at which  $timer_j$  is set to such a value. Then, for any  $\tau \geq \tau_1$  and any  $x \geq x_0$  we have that  $T_R(\tau, x) \geq f_R(\tau, x) \geq f_R(\tau_f, x_0) = \Delta'$ , from assumptions (f3) and (f1), and the above paragraph, respectively.

Let us now consider any two consecutive executions of the lines 15 to 25 by  $p_j$  with  $k = i$  that start at time  $\tau > \max(\tau_{01}, \tau_1)$ . Let us notice that, by the definition of  $\tau$ , these two executions are separated by at least  $\Delta'$  time units. We claim that in the second such execution  $SUSPICIONS[j, i]$  is not increased. Since this holds for any two executions and any time  $\tau$ , the lemma follows. We consider two cases:

- In the second execution, the value of  $STOP[i]$  that is read by  $p_j$  (line 15) is *true*. In that case, independently of the result of the test at line 17,  $SUSPICIONS[j, i]$  will not be incremented in that execution, because of the test at line 20.
- In the second execution, the value of  $STOP[i]$  that is read by  $p_j$  (line 15) is *false*. When this read operation takes effect, task T2 of process  $p_i$  was either between two updates of  $PROGRESS[i]$  (line 8) adjacent in  $S$ , or between an update of  $PROGRESS[i]$  and an update of  $STOP[i]$  at line 11, both adjacent in  $S$ . In either case, the value of  $PROGRESS[i]$  was updated at most  $3\Delta$  time units before the read of  $STOP[i]$  by  $p_j$  takes effect. Since the latest update of  $last_j[i]$  was done at least  $\Delta'$  time units before, and  $\Delta' > 3\Delta$ , the value read in  $progress_k[i]$  and the value in  $last_j[i]$  are different. Then, the test at line 17 evaluates to true and  $SUSPICIONS[j, i]$  is not incremented.

<sup>9</sup>Counting  $\Delta$  for the first execution of line 8,  $\Delta$  for the possible accesses of  $STOP[i]$  at line 9, and  $\Delta$  for the second execution of line 8. Let us notice that we do not count the time of a read operation at these lines. As indicated in Section 3.2, this is because, as the critical variables  $PROGRESS[i]$  and  $STOP[i]$  are written only by  $p_i$ , that process can manage local copies and read these local copies instead of reading the corresponding critical variables. The important point here is the determination of an upper bound for the time that can elapse between the completion of two consecutive write operations on the critical variables of  $p_i$ . The same remark applies to the computation of the next time duration, where  $\Delta$  time units are counted for each of the lines 8, 9, and 11.

Let  $(M_\ell, \ell) = \text{lexmin}(\{M_x, x \mid x \in B\})$ .

**Lemma 3** *There is a single process  $p_\ell$  and it is correct.*

**Proof** The lemma follows directly from the following observations:  $B$  does not contain faulty processes (definition),  $B \neq \emptyset$  (Lemma 2), and no two processes have the same identity (initial assumption). □ Lemma 3

**Lemma 4** *There is a time after which  $p_\ell$  permanently executes the loop defined by the lines 7-10 of task T2.*

**Proof** Due to Lemma 1, for each faulty process  $p_k$ , there is a time after which the predicate  $k \notin \text{candidate}_\ell$  remains forever true. So, after that time the faulty processes do not compete with  $p_\ell$  to become leader (line 4). This constitutes observation O1.

Let us now consider a correct process  $p_i$ . It follows from the definition of  $B$  that we have the following: (1) if  $i \in B$ , eventually  $\sum_{1 \leq j \leq n} \text{SUSPICIONS}[j, i] = M_i \geq M_\ell$ ; (2) if  $i \notin B$ , eventually  $\sum_{1 \leq j \leq n} \text{SUSPICIONS}[j, i] > M_\ell$ . This constitutes observation O2.

It follows from the previous observations O1 and O2, and the fact that we always have  $x \in \text{candidate}_x$  for any process  $p_x$ , that after some finite time, each evaluation of the predicate  $\text{leader}() = \ell$  (line 7), returns true to  $p_\ell$ . Consequently, there is a time after which  $p_\ell$  permanently executes the lines 7-10 of task T2. □ Lemma 4

**Theorem 1** *There is a time after which a correct process is elected as the eventual common leader.*

**Proof** We show that  $p_\ell$  is the eventual common leader. From Lemma 3  $p_\ell$  is unique and correct. Moreover, due to the definitions of the bound  $M_\ell$  and the set  $B$ , there is a finite time  $\tau$  after which, for each correct process  $p_i$ ,  $i \neq \ell$ , we have  $(\sum_{1 \leq j \leq n} \text{SUSPICIONS}[j, i], i) > (M_\ell, \ell)$ . Moreover, due to Lemma 1, there is a time after which, for each correct process  $p_i$  and each faulty process  $p_k$  we have  $k \notin \text{candidate}_i$ . It follows from these observations, that proving the theorem amounts to show that eventually the predicate  $\ell \in \text{candidate}_i$  remains permanently true at each correct process  $p_i$ .

Let us notice that the predicate  $x \in \text{candidate}_x$  is always true for any process  $p_x$ . This follows from the fact that initially  $x$  belongs to  $\text{candidate}_x$ , and then  $p_x$  does not execute the tasks T3 for  $k = x$ , and consequently cannot withdraw  $x$  from  $\text{candidate}_x$ . It follows that we always have  $\ell \in \text{candidate}_\ell$ . So, let us examine the case  $i \neq \ell$ .

It follows from Lemma 4 that there is a time  $\tau$  after which  $p_\ell$  remains permanently in the **while** loop of task T2. Let  $\tau' \geq \tau$  be a time at which we have  $\sum_{1 \leq j \leq n} \text{SUSPICIONS}[j, \ell] = M_\ell$ , and  $p_\ell$  has executed line 9 (i.e.,  $\text{STOP}[\ell]$  remains false forever).

After  $\tau'$ , because  $p_\ell$  is forever increasing  $\text{PROGRESS}[\ell]$ , the test of line 17 eventually evaluates to true and (if not already done)  $p_i$  adds  $\ell$  to  $\text{candidate}_i$ . We claim that, after that time, the task T3 of  $p_i$  is always executing the lines 18-19, from which it follows that  $\ell$  remains forever in  $\text{candidate}_i$ .

*Proof of the claim.* Let us assume by contradiction that the test of line 17 is false when evaluated by  $p_i$ . It follows that  $\ell$  is withdrawn from  $\text{candidate}_i$ , and this occurs at line 24. (It cannot occur at line 21 because after  $\tau$  we always have  $\text{STOP}[\ell] = \text{false}$ .) But line 23 is executed before 24, from which we conclude that  $\text{SUSPICIONS}[i, \ell]$  has been increased, which means that we have now  $\sum_{1 \leq j \leq n} \text{SUSPICIONS}[j, \ell] = M_\ell + 1$ , contradicting the definition of the bound  $M_\ell$ . *End of the proof of the claim.* □ Theorem 1

**Theorem 2** *Let  $p_\ell$  be the eventual common leader. All shared variables (but  $\text{PROGRESS}[\ell]$ ) are bounded.*

**Proof** Let us consider a time  $\tau$  after which the eventual common leader has been elected. Due to Theorem 1, such a time does exist. After  $\tau$ , as it is not the leader, no process  $p_k$ ,  $k \neq \ell$ , can execute the **while** loop of task T2, and consequently  $\text{PROGRESS}[k]$  is no longer increased.

It follows from the fact that  $\text{PROGRESS}[k]$  is no longer increased that, after some time, for any  $p_j$ , we have  $\text{last}_j[k] = \text{PROGRESS}[k]$ . Consequently, after that time, the test of line 17 is never satisfied and, if  $k$  belongs to  $\text{candidates}_j$ , it is suppressed from this set and never added again. It follows that, as far as  $\text{SUSPICIONS}[j, k]$  is concerned, line 23 is no longer executed, which means that  $\text{SUSPICIONS}[j, k]$  is no longer increased.

The fact that, for any  $j$ ,  $\text{SUSPICIONS}[j, \ell]$  is bounded has been proved in Lemma 2. □ Theorem 2

**Theorem 3** *After a finite time, only one process (the eventual common leader) writes forever into the shared memory. Moreover, it always writes the same shared variable.*

**Proof** After an eventual common leader  $p_\ell$  had been elected, that leader executes permanently the **while** loop of task  $T_2$  and consequently forever increases  $PROGRESS[\ell]$ . After it has entered the loop it sets  $STOP[\ell]$  to *false*, and we can conclude from line 9 that it will not longer write  $STOP[\ell]$ .

The fact that, after some time, no variable  $SUSPICIONS[j, k]$  ( $1 \leq j, k \leq n$ ) is written, follows directly from the combined effect of (1) the fact that each of these variables is bounded (Theorem 2), and (2) each write increases its value. The same reasoning applies to each variable  $PROGRESS[k]$  for  $k \neq \ell$ .

Let us now consider a  $STOP[k]$  variable,  $k \neq \ell$ . As, after the eventual common leader  $p_\ell$  has been elected,  $p_k$  no longer executes the body of the **while** loop of task  $T_2$ , it follows that  $STOP[k]$  cannot be written at line 9. Let us now consider the write at line 11. If  $STOP[k]$  is equal to *false*,  $p_k$  sets it to *true*. Then task  $T_2$  of  $p_k$  always finds  $STOP[k]$  equal to *true* at line 11, and consequently never updates it again. □ *Theorem 3*

### 3.4 Optimality Results

Let  $\mathcal{A}$  be any algorithm that implements  $\Omega$  in  $\mathcal{AS}_n[AWB]$  with up to  $t$  faulty processes. We have the following lower bounds.

**Lemma 5** *Let  $R$  be any run of  $\mathcal{A}$  with less than  $t$  faulty processes and let  $p_\ell$  be the leader chosen in  $R$ . Then  $p_\ell$  must write forever in the shared memory in  $R$ .*

**Proof** Assume, by way of contradiction, that  $p_\ell$  stops writing in the shared memory in run  $R$  at time  $\tau$ . Consider another run  $R'$  of  $\mathcal{A}$  in which all processes behave like in  $R$  except  $p_\ell$ , which behaves exactly like in  $R$  until time  $\tau + 1$ , and crashes at that time. Since at most  $t$  processes crash in  $R'$ , by definition of  $\mathcal{A}$ , eventually a leader must be elected. In fact, in  $R'$  all the processes except  $p_\ell$  behave exactly like in  $R$  and elect  $p_\ell$  as their (permanent) leader. These processes cannot distinguish  $R'$  from  $R$  and cannot detect the crash of  $p_\ell$ . Hence, in  $R'$  algorithm  $\mathcal{A}$  does not satisfy the Eventual Leadership property of  $\Omega$ , which is a contradiction. Therefore,  $p_\ell$  cannot stop writing in the shared memory. □ *Lemma 5*

**Lemma 6** *Let  $R$  be any run of  $\mathcal{A}$  with less than  $t$  faulty processes and let  $p_\ell$  be the leader chosen in  $R$ . Then every correct process  $p_i$ ,  $i \neq \ell$ , must read forever from the shared memory in  $R$ .*

**Proof** Assume, by way of contradiction, that a correct process  $p_i$  stops reading from the shared memory in run  $R$  at time  $\tau$ . Let  $\tau'$  be the time at which  $p_i$  chooses permanently  $p_\ell$  as leader. Consider another run  $R'$  of  $\mathcal{A}$  in which  $p_\ell$  behaves exactly like in  $R$  until time  $\max(\tau, \tau') + 1$ , and crashes at that time. Since at most  $t$  processes crash in  $R'$ , by definition of  $\mathcal{A}$ , a leader must be eventually elected. In  $R'$ , we make  $p_i$  to behave exactly like in  $R$ . As it stopped reading the shared memory at time  $\tau$ ,  $p_i$  cannot distinguish  $R'$  from  $R$  and cannot detect the crash of  $p_\ell$ . Hence in  $R'$ ,  $p_i$  elects  $p_\ell$  as its (permanent) leader at time  $\tau'$ . Hence, in  $R'$  algorithm  $\mathcal{A}$  does not satisfy the Eventual Leadership property of  $\Omega$ , which is a contradiction. Therefore,  $p_i$  cannot stop reading from the shared memory. □ *Lemma 6*

The following theorem follows immediately from the previous lemmas.

**Theorem 4** *The algorithm described in Figure 2 is optimal in with respect to the number of processes that have to write the shared memory. It is quasi-optimal with respect to the number of processes that have to read the shared memory.*

The “quasi-optimality” comes from the fact that the algorithm described in Figure 2 requires that each process (including the leader) reads forever the shared memory (all the processes have to read the array  $SUSPICIONS[1..n, 1..n]$ ).

### 3.5 Discussion

**Using multi-writer/multi-reader ( $nWnR$ ) atomic registers** If we allow  $nWnR$  atomic variables, each column  $SUSPICIONS[*, j]$  can be replaced by a single  $SUSPICIONS[j]$ . Consequently vectors of  $nWnR$  atomic variables can be used instead of matrices of  $1WnR$  atomic variables.

**Eliminating the local clocks** The timers (and consequently the local clocks used to implement them) can be eliminated if we consider that each execution of the statement  $timer_i \leftarrow timer_i - 1$  takes at least one time unit. The code of task  $T3$  becomes then the following:

```

task T3: timer_i ← 1;
  while (true) do timer_i ← timer_i - 1;
    if (timer_i = 0) then Line 14 until Line 26 of Figure 2 or 5;
      timer_i ← max({SUSPICIONS[i, k]}_{1 ≤ k ≤ n}) + 1
    end_if
  end_while.

```

## 4 An $\Omega$ algorithm for $\mathcal{AS}_n[AWB]$ with Bounded Variables Only

### 4.1 A Lower Bound Result

This section shows that any algorithm that implements  $\Omega$  in  $\mathcal{AS}_n[AWB]$  with only bounded memory requires all correct processes to read and write the shared memory forever. As we will see, it follows from this lower bound that the algorithm described in Figure 5 is optimal with respect to this criterium.

Let  $\mathcal{A}$  be an algorithm that implements  $\Omega$  in  $\mathcal{AS}_n[AWB]$  such that, in every run  $R$  of  $\mathcal{A}$ , the number of shared memory bits used is bounded by a value  $S_R$  (which may depend on the run). This means that in any run there is time after which no new memory positions are used, and each memory position has bounded number of bits. To make the result stronger, we also assume that  $\mathcal{A}$  knows  $t$  (maximum number of processes that can fail in any run of  $\mathcal{A}$ ).

**Theorem 5** *The algorithm  $\mathcal{A}$  has runs in which at least  $t + 1$  processes write forever in the shared memory.*

**Proof** To prove the claim we construct a run  $R$  of  $\mathcal{A}$  such that:

1.  $R$  is fault free,
2. Process  $p_1$  is synchronous while the rest of processes are asynchronous, and
3. There is an infinite sequence of times  $\tau_0 < \tau_1 < \tau_2 < \dots$  such that,  $\forall i > 0$ , in the interval  $(\tau_{i-1}, \tau_i]$  some process changes its leader or at least  $t + 1$  processes write in the shared memory.

Clearly, since a leader must be eventually elected in  $R$  and the number of processes is finite, due to Item 3, there is a set of at least  $t + 1$  processes that write in the shared memory forever.

For simplicity, let us define  $\tau_0 = 0$ . This will be the base case. Then, for  $i > 0$  let us assume  $R$  is already constructed up to time  $\tau_{i-1}$ . We construct now interval  $(\tau_{i-1}, \tau_i]$ . This interval is constructed differently depending on which of the following two cases occurs.

- If at time  $\tau_{i-1}$  the leader of some process  $p_j$  is an asynchronous process  $p_k$  (i.e.,  $k \neq 1$ ), we first consider a run  $R_i$  that behaves exactly like  $R$  up to time  $\tau_{i-1}$ . Then, after that time all processes advance synchronously (e.g., one step per time unit), except  $p_k$  which crashes at time  $\tau_{i-1} + 1$ . By Eventual Leadership, there is a time  $\tau > \tau_{i-1}$  in  $R_i$  at which no process considers  $p_k$  as its leader. Then, let us define  $\tau_i = \tau + 1$  and make  $R$  to behave in the interval  $(\tau_{i-1}, \tau_i]$  as follows. All processes except  $p_k$  behave in this interval exactly like in the interval  $(\tau_{i-1}, \tau_i]$  of  $R_i$ . Process  $p_k$  does not crash, but is stopped at time  $\tau_{i-1} + 1$  and does not execute any step until the end of the interval. This behavior is possible since  $p_k$  is asynchronous. Then, we have that in the interval  $(\tau_{i-1}, \tau_i]$  some process changed its leader. This ends the first case.
- The second case occurs when at time  $\tau_{i-1}$  in  $R$  the leader of all processes is the synchronous process  $p_1$ . As before we now consider an auxiliary run  $R_i$  that behaves exactly like  $R$  up to time  $\tau_{i-1}$ . After that time all processes advance synchronously (e.g., one step per time unit) in  $R_i$ . If some process  $p_j$  changes its leader in  $R_i$  at some time  $\tau > \tau_{i-1}$ , then we define  $\tau_i = \tau + 1$  and make the interval  $(\tau_{i-1}, \tau_i]$  of  $R$  behave exactly as interval  $(\tau_{i-1}, \tau_i]$  of  $R_i$ .

Otherwise, if no process changes its leader in  $R_i$  after  $\tau_{i-1}$ , we have from Lemma 5 that  $p_1$  writes in the shared memory forever. Let us assume by way of contradiction that there is a time  $\tau > \tau_{i-1}$  after which at most  $t - 1$  other processes write forever in the shared memory in  $R_i$ . Since the shared memory is bounded, some state

(understood as the value of all its bits)  $S$  of the shared memory must occur infinitely often in  $R_i$  after  $\tau$ . (First line in Figure 4 where the state  $S$  is described with an area with stripes.)

Let us consider now a run  $R'_i$  which behaves exactly like  $R_i$  up to time  $\tau' > \tau$  at which the shared memory is in state  $S$  (second line in Figure 4). Then, at that time the (up to  $t$ ) processes that were writing in the shared memory (including  $p_1$ ) crash in  $R'_i$ . The rest of the processes advance synchronously (and hence the  $AWB_1$  assumption holds in  $R'_i$ ) until the smallest time  $\tau'' > \tau'$  at which some process changes its leader or some process writes in the shared memory. This must eventually occur by Eventual Leadership, since the leader of all the processes at time  $\tau'$  has crashed in  $R'_i$ . Note that in the interval  $(\tau', \tau'')$  all read operations find the shared memory in state  $S$ .

Consider now another run  $R''_i$  in which the up to  $t$  processes (including  $p_1$ ) that write forever in  $R_i$  behave like they do in that run, while the rest of processes (let us denote this set of processes by  $L$ ) behave like in  $R_i$  up to time  $\tau'$  (last line in Figure 4.) After  $\tau'$ , the processes in  $L$  are delayed (note that they are all asynchronous) so that every time they read from the shared memory they find it in state  $S$  (see Figure 4). From the behavior of the processes in  $L$  in run  $R'_i$  and the fact that they cannot distinguish run  $R''_i$  from run  $R'_i$ , we have that there is a time  $\tau''' > \tau'$  at which some process in  $L$  changes its leader or writes in the shared memory in run  $R''_i$ . Then, we define  $\tau_i = \tau''' + 1$  and make interval  $(\tau_{i-1}, \tau_i]$  of  $R$  behave exactly like that interval in  $R''_i$ .

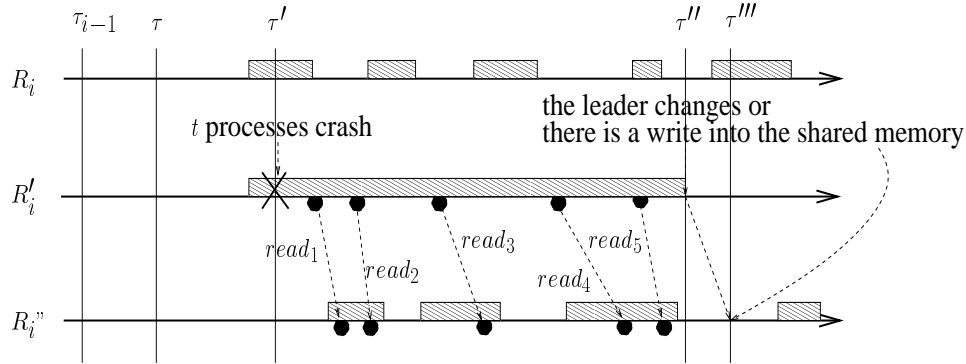


Figure 4: Illustrating the runs  $R_i$ ,  $R'_i$  and  $R''_i$

Figure 4 summarizes the previous reasoning. In the first run  $R_i$ , after  $\tau$ , only  $t$  processes write forever. The same state  $S$  (depicted by the area with stripes) occurs repeatedly forever. In the run  $R'_i$ , these  $t$  processes crash in state  $S$  (they crash at the time marked with a cross). The read operations from the other processes are indicated with black dots. In the run  $R''_i$ , the same processes as in  $R'_i$  read while the system in the state  $S$ .

□*Theorem 5*

The system model defined in this paper assumes  $t = n - 1$ . Hence the following corollary.

**Corollary 1** *Any algorithm that implements  $\Omega$  in  $\mathcal{AS}_n[AWB]$  with bounded shared memory has runs in which all processes write the shared memory forever.*

## 4.2 An algorithm with only bounded variables

**Principles and description** As already indicated, we are interested here in an algorithm whose variables are all bounded. To attain this goal, we use a hand-shaking mechanism. More precisely, we replace the shared array  $PROGRESS[1..n]$  and all the local arrays  $last_i[1..n]$ ,  $1 \leq i \leq n$ , by two shared matrices of  $1WnR$  boolean values, denoted  $PROGRESS[1..n, 1..n]$  and  $LAST[1..n, 1..n]$ .

The hand-shaking mechanism works as follows. Given a pair of processes  $p_i$  and  $p_k$ ,  $PROGRESS[i, k]$  and  $LAST[i, k]$  are used by these processes to send signals to each other. More precisely, to signal  $p_k$  that it is alive,  $p_i$  sets  $PROGRESS[i, k]$  equal to  $\neg LAST[i, k]$ . In the other direction,  $p_k$  indicates that it has seen this “signal” by cancelling it, namely, it resets  $LAST[i, k]$  equal to  $PROGRESS[i, k]$ . It follows from the essence of the hand-shaking

mechanism that both  $p_i$  and  $p_k$  have to write shared variables, but as shown by Corollary 1, this is the price that has to be paid to have bounded shared variables.

Using this simple technique, we obtain the algorithm described in Figure 5. In order to capture easily the parts that are new or modified with respect to the previous algorithm, the line number of the new statements are suffixed with the letter R (so the line 8 of the previous protocol is replaced by three new lines, while each of the lines 16, 17 and 19 is replaced by a single line). This allows a better understanding of the common principles on which both algorithms rely.

```

task T1:
(1) when leader() is invoked:
(2)   for_each  $k \in candidates_i$  do
(3)      $susp_i[k] \leftarrow \sum_{1 \leq j \leq n} SUSPICIONS[j, k]$  end_for;
(4)   let  $(-, \ell) = \text{lex\_min}(\{(susp_i[k], k)\}_{k \in candidates_i});$ 
(5)   return( $\ell$ )

task T2:
(6) repeat_forever
(7)   while (leader() =  $i$ ) do
(8.R1)     for_each  $k \in \{1, \dots, n\} \setminus \{i\}$  do
(8.R2)        $PROGRESS[i, k] \leftarrow \neg LAST[i, k]$ 
(8.R3)     end_for;
(9)     if  $STOP[i]$  then  $STOP[i] \leftarrow false$  end_if
(10)    end_while;
(11)    if  $(\neg STOP[i])$  then  $STOP[i] \leftarrow true$  end_if
(12)  end_repeat

task T3:
(13) when  $timer_i$  expires:
(14)   for_each  $k \in \{1, \dots, n\} \setminus \{i\}$  do
(15)      $stop\_k_i \leftarrow STOP[k];$ 
(16.R1)   $progress\_k_i \leftarrow PROGRESS[k, i];$ 
(17.R1)  if  $(progress\_k_i \neq LAST[k, i])$  then
(18)      $candidates_i \leftarrow candidates_i \cup \{k\};$ 
(19.R1)   $LAST[k, i] \leftarrow progress\_k_i$ 
(20)   else_if  $(stop\_k_i)$  then
(21)      $candidates_i \leftarrow candidates_i \setminus \{k\}$ 
(22)   else_if  $(k \in candidates_i)$  then
(23)      $SUSPICIONS[i, k] \leftarrow SUSPICIONS[i, k] + 1;$ 
(24)      $candidates_i \leftarrow candidates_i \setminus \{k\}$ 
(25)   end_if
(26)  end_for;
(27)  set  $timer_i$  to  $\max(\{SUSPICIONS[i, k]\}_{1 \leq k \leq n}) + 1$ 

```

Figure 5: All variables are 1WMR and bounded

**Proof of the algorithm** The statement of the lemmas 1, 2, 3 and 4, and Theorem 1 are still valid when the shared array  $PROGRESS[1..n]$  and the local arrays  $last_i[1..n]$ ,  $1 \leq i \leq n$  are replaced by the shared matrices  $PROGRESS[1..n, 1..n]$  and  $LAST[1..n, 1..n]$ .

As far as their proofs are concerned, the proofs of the lemmas 3 and 4 given in Section 3.3 are verbatim the same. The proofs of the lemmas 1 and 2, and the proof of Theorem 1 have to be slightly modified to suit to the new context. Basically, they differ from their counterparts of Section 3.3 in the way they establish the property that, after some time, no correct process  $p_i$  misses an “alive” signal from a process that satisfies the assumption  $AWB_1$ . (More specifically, the sentence “there is a time after which  $PROGRESS[k]$  does no longer increase” has to be replaced by the sentence “there is a time after which  $PROGRESS[k, i]$  remains forever equal to  $LAST[k, i]$ ”). As they are very close to the previous ones and tedious, we don’t detail these proofs. (According to the usual sentence, “They are left as an exercise to the reader”.)

The same reasoning as the one done in the proof of the Theorem 2 shows that each shared variable  $SUSPICIONS[j, k]$ ,  $1 \leq j, k \leq n$ , is bounded. Combined with the fact that the variables  $PROGRESS[j, k]$  and  $LAST[j, k]$  are boolean, we obtain the following theorem.

**Theorem 6** *All the variables used in the algorithm described in Figure 5 are bounded.*

The following theorem is the counterpart of Theorem 3.

**Theorem 7** *Let  $p_\ell$  be the process elected as the eventual common leader, and  $p_i, i \neq \ell$ , any correct process. There is a time after which the only variables that are written are  $PROGRESS[\ell, i]$  (written by  $p_\ell$ ) and  $LAST[\ell, i]$  (written by  $p_i$ ).*

**Proof** The proof that the variables  $PROGRESS[\ell, j]$ ,  $1 \leq j \leq n$ , are infinitely often written, and the proof that there is a time after which the variables  $STOP[j]$ ,  $1 \leq j \leq n$ , and the variables  $SUSPICIONS[j, k]$ ,  $1 \leq j, k \leq n$ , are no longer written is the same as the proof done in Theorem 3.

The fact that there is a time after which  $PROGRESS[x, j]$ ,  $1 \leq x, j \leq n, x \neq \ell$ , are no longer written follows from the fact that, after  $p_\ell$  has been elected, no process  $p_x$  executes the body of the **while** loop of task  $T2$ .

Let us now consider any variable  $LAST[x, y]$ ,  $x \neq \ell$ . As, after  $p_\ell$  has been elected, no correct process  $p_x, x \neq \ell$ , updates  $PROGRESS[x, y]$  (at line 8.R2), it follows that there is a time after which  $LAST[x, y] = PROGRESS[x, y]$  remains forever true for  $1 \leq x, y \leq n$  and  $x \neq \ell$ . Consequently, after a finite time, the test of line 17.R1 is always false for  $p_x, x \neq \ell$ , and  $LAST[x, y]$  is no longer written.  $\square_{Theorem 7}$

Finally, the next theorem follows directly from Corollary 1.

**Theorem 8** *The  $\Omega$  algorithm described in Figure 5 is optimal with respect to the number of processes that have to write the shared memory.*

## 5 Conclusion

This paper has addressed the problem of electing an eventual leader in an asynchronous shared memory system. It has three main contributions.

- The first contribution is the statement of an assumption (a property denoted  $AWB$ ) that allows electing a leader in the shared memory asynchronous systems that satisfy that assumption. This assumption requires that after some time (1) there is a process whose write accesses to some shared variables are timely, and (2) the other processes have asymptotically well-behaved timers. The notion of asymptotically well-behaved timer is weaker than the usual notion of timer where the timer durations have to monotonically increase when the values to which they are set increase. This means that  $AWB$  is a particular weak assumption.
- The second contribution is the design of two algorithms that elect an eventual leader in any asynchronous shared memory system that satisfies the assumption  $AWB$ . In addition of being independent of  $t$  (the maximum number of processes allowed to crash), and being based only on one-writer/multi-readers atomic shared variables, these algorithms enjoy noteworthy properties. The first algorithm guarantees that (1) there is a (finite) time after which a single process writes forever the shared memory, and (2) all but one shared variables have a bounded domain. The second algorithm uses (1) a bounded memory but (2) requires that each process forever writes the shared memory.
- The third contribution shows that the previous tradeoff (bounded/unbounded memory vs number of processes that have to write) is inherent to the leader election problem in asynchronous shared memory systems equipped with  $AWB$ . It follows that both algorithms are optimal, the first with respect to the number of processes that have to forever write the shared memory, the second with respect to the boundedness of the memory.

Several questions remain open. One concerns the first algorithm. Is it possible to design a leader algorithm in which there is a time after which the eventual leader is not required to read the shared memory? Another question is the following: is the second algorithm optimal with respect to the size of the control information (bit arrays) it uses to have a bounded memory implementation?

## References

- [1] Abraham I., Chockler G.V., Keidar I. and Malkhi D., Byzantine Disk Paxos, Optimal Resilience with Byzantine Shared Memory. *Proc. 23th ACM Symposium on Principles of Distributed Computing (PODC'04)*, ACM Press, pp. 226-235, 2004.
- [2] Aguilera M.K., Delporte-Gallet C., Fauconnier H. and Toueg S., On Implementing Omega with Weak Reliability and Synchrony Assumptions. *Proc. 22th ACM Symposium on Principles of Distributed Computing (PODC'03)*, ACM Press, pp. 306-314, 2003.
- [3] Aguilera M.K., Delporte-Gallet C., Fauconnier H. and Toueg S., Communication-Efficient Leader Election and Consensus with Limited Link Synchrony. *Proc. 23th ACM Symposium on Principles of Distributed Computing (PODC'04)*, ACM Press, pp. 328-337, 2004.
- [4] Aguilera M.K., Englert B. and Gafni E., On Using Network Attached Disks as Shared Memory. *Proc. 21th ACM Symposium on Principles of Distributed Computing (PODC'03)*, ACM Press, pp. 315-324, 2003.
- [5] Chandra T. and Toueg S., unreliable Failure Detectors for Resilient Distributed Systems. *Journal of the ACM*, 43(2):225-267, 1996.
- [6] Chandra T., Hadzilacos V. and Toueg S., The Weakest Failure Detector for Solving Consensus. *Journal of the ACM*, 43(4):685-722, 1996.
- [7] Dwork C., Lynch N. and Stockmeyer L., Consensus in the Presence of Partial Synchrony. *Journal of the ACM*, 35(2):288-323, 1988.
- [8] Fischer M.J., Lynch N. and Paterson M.S., Impossibility of Distributed Consensus with One Faulty Process. *Journal of the ACM*, 32(2):374-382, 1985.
- [9] Gafni E. and Lamport L., Disk Paxos. *Distributed Computing*, 16(1):1-20, 2003.
- [10] Gibson G.A. *et al.*, A Cost-effective High-bandwidth Storage Architecture. *Proc. 8th Int'l Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS'98)*, ACM Press, pp. 92-103, 1998.
- [11] Guerraoui R. and Raynal M., The Information Structure of Indulgent Consensus. *IEEE Transactions on Computers*, 53(4):453-466, 2004.
- [12] Guerraoui R. and Raynal M., The Alpha of Asynchronous Consensus. *The Computer Journal*, To appear, 2006.
- [13] Guerraoui R. and Raynal M., A Leader Election Protocol for Eventually Synchronous Shared Memory Systems. *4th Int'l IEEE Workshop on Software Technologies for Future Embedded and Ubiquitous Systems (SEUS'06)*, IEEE Computer Society Press, pp. 75-80, 2006.
- [14] H elary J.-M., Hurfin M., Mostefaoui A., Raynal M. and Tronel F., Computing Global Functions in Asynchronous Distributed Systems with Perfect Failure Detectors. *IEEE Transactions on Parallel and Distributed Systems*, 11(9):897-909, 2000.
- [15] Herlihy M.P. and Wing J.M., Linearizability: a Correctness Condition for Concurrent Objects. *ACM Transactions on Programming Languages and Systems*, 12(3):463-492, 1990.
- [16] Lamport L., The Part-Time Parliament. *ACM Transactions on Computer Systems*, 16(2):133-169, 1998. (The first version of Paxos appeared a a DEC Tech Report in 1989.)
- [17] Larrea M., Fern andez A. and Ar evalo S., Optimal Implementation of the Weakest Failure Detector for Solving Consensus. *Proc. 19th Symposium on Resilient Distributed Systems (SRDS'00)*, IEEE Computer Society Press, pp. 52-60, 2000.
- [18] Lee E.K. and Thekkath C., Petal: Distributed Virtual Disks. *Proc. 7th Int'l Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS'96)*, ACM Press, pp. 84-92, 1996.
- [19] Lo W.-K. and Hadzilacos V., Using failure Detectors to solve Consensus in Asynchronous Shared Memory Systems. *Proc. 8th Int'l Workshop on Distributed Computing (WDAG'94)*, Springer Verlag LNCS #857, pp. 280-295, 1994.
- [20] Malkhi D., Oprea F. and Zhou L.,  $\Omega$  Meets Paxos: Leader Election and Stability without Eventual Timley Links. *Proc. 19th Int'l Symposium on Distributed Computing (DISC'05)*, Springer Verlag LNCS #3724, pp. 199-213, 2005.
- [21] Mostefaoui A., Mourgaya E., and Raynal M., Asynchronous Implementation of Failure Detectors. *Proc. Int'l IEEE Conference on Dependable Systems and Networks (DSN'03)*, IEEE Computer Society Press, pp. 351-360, 2003.

- [22] Mostefaoui A. and Raynal M., Leader-Based Consensus. *Parallel Processing Letters*, 11(1):95-107, 2001.
- [23] Mostéfaoui A., Raynal M. and Travers C., Crash Resilient Time-Free Eventual Leadership. *Proc. 23th Symposium on Resilient Distributed Systems (SRDS'04)*, IEEE Computer Society Press, pp. 208-218, 2004.
- [24] Mostéfaoui A., Raynal M. and Travers C., Time-free and Timeliness Assumptions can be Combined to Get Eventual Leadership. *IEEE Transactions on Parallel and Distributed Systems*, 17(7):656-666, 2006.
- [25] Powell D., Failure Mode Assumptions and Assumption Coverage. *Proc. of the 22nd Int'l Symposium on Fault-Tolerant Computing (FTCS-22)*, Boston, MA, pp.386-395, 1992.
- [26] Raynal M., A Short Introduction to Failure Detectors for Asynchronous Distributed Systems. *ACM SIGACT News, Distributed Computing Column*, 36(1):53-70, 2005.